



Middlebury

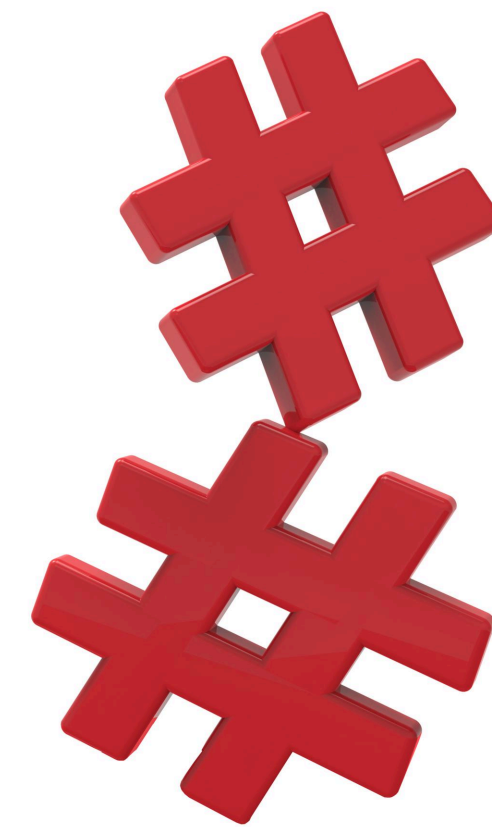
CSCI 201: Data Structures

Spring 2025

Lecture 10M: Hash Tables (Part 1)

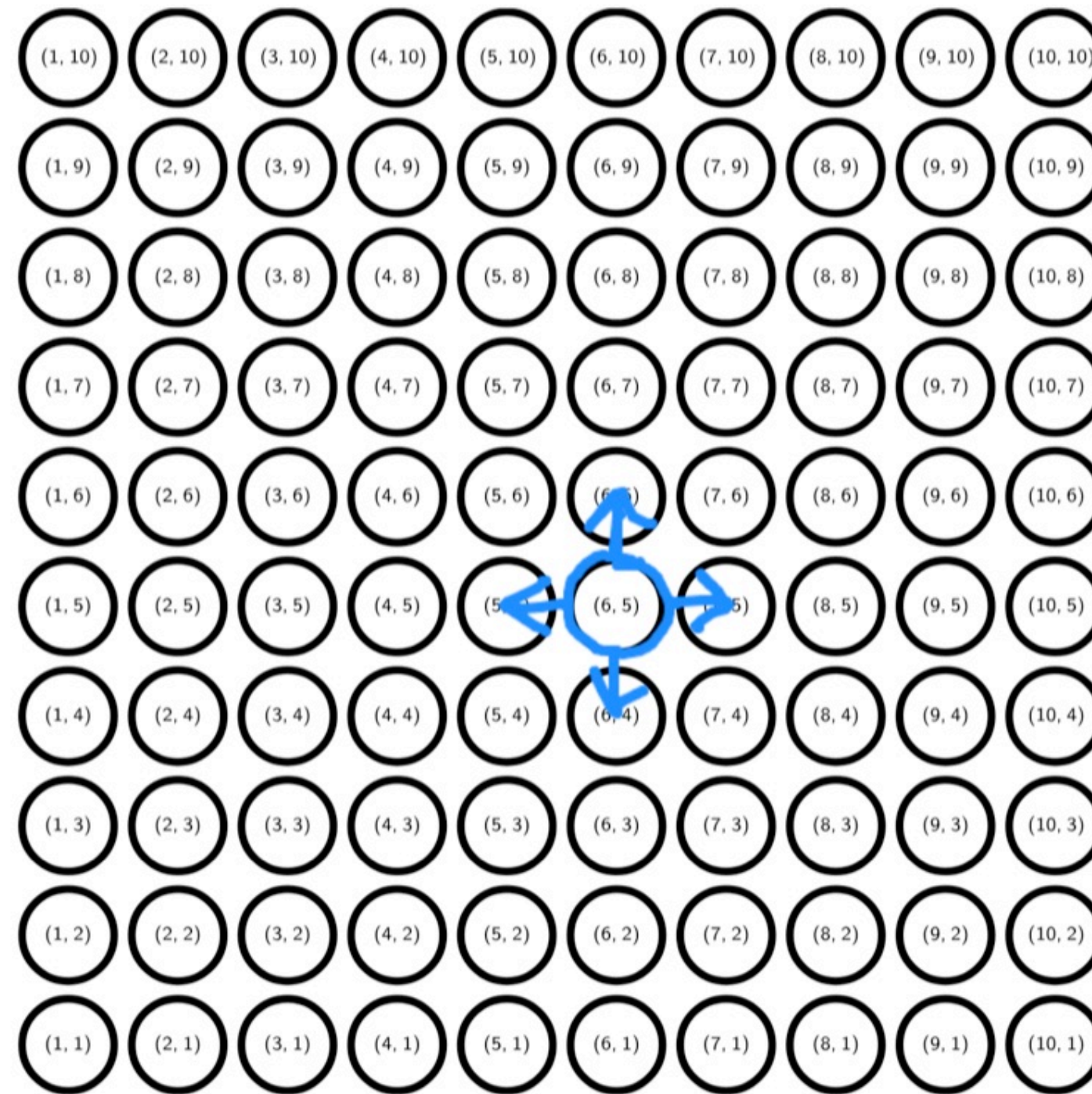
Goals for today:

- Solve problems using **Java**'s built-in **HashSet** and **HashMap** (which are built using a **hash table**).
- Write our own **hashCode** and **equals** methods for our custom types in order to use them with **Java**'s built-in **HashSet** and **HashMap**.
- Use a hash function to determine the index of a key in a hash table.
- Practice with bitwise operators.



Consider this problem:

- Suppose we start at some point in a grid and can move left, right, up or down (no diagonal movement).
- We can only step into a grid point that we have not yet visited and cannot step out of bounds of the grid.
- Open `PointTracker.java` and brainstorm how you would design a solution for the `boolean haveBeenHere(Point point)` method (without a `Set` or `Map`).



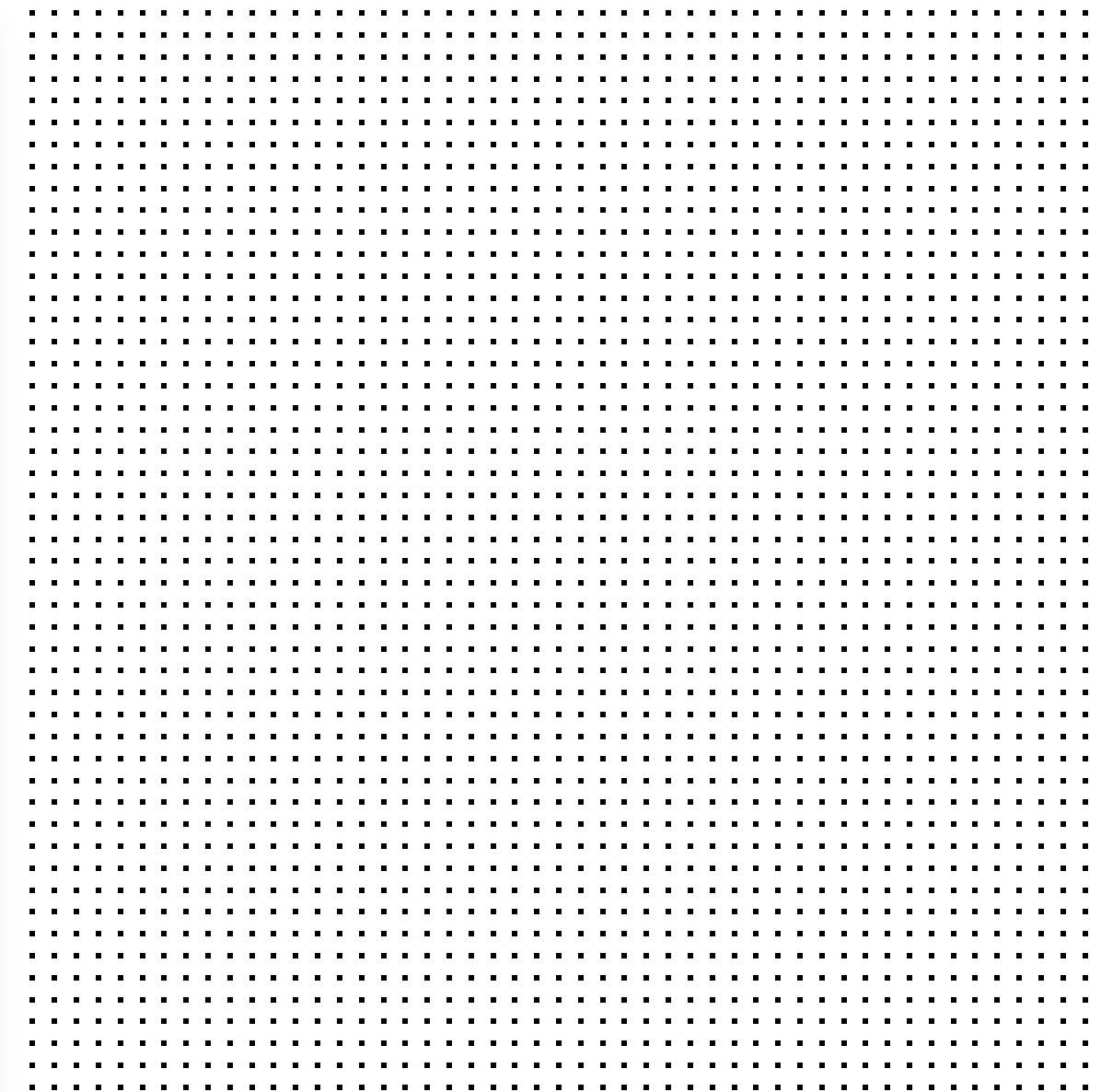
Each iteration :

- try random direction
- out of bounds?
- visited here already?
- stuck?

One possible solution:

but...

```
1 public class PointTracker {
2     boolean[][] visited;
3
4     public PointTracker() {
5         visited = new boolean[GRID_SIZE][GRID_SIZE];
6         for (int i = 0; i < GRID_SIZE; i++) {
7             for (int j = 0; j < GRID_SIZE; j++) {
8                 visited[i][j] = false;
9             }
10        }
11        location = new Point(0, 0); // or in middle
12        move(location);
13    }
14
15    public boolean haveBeenHere(Point point) {
16        return visited[point.x][point.y];
17    }
18
19    public void move(Point newLocation) {
20        visited[newLocation.x][newLocation.y] = true;
21        location = newLocation;
22    }
23 }
```



Is there a way to use less memory?

We only want to keep track of points we've visited.

ArrayList of **Points**?



- add new location $O(1)$
- visited before? $O(\# \text{ of points visited})$

But we also want to be able to (quickly) look up if we've visited a point.

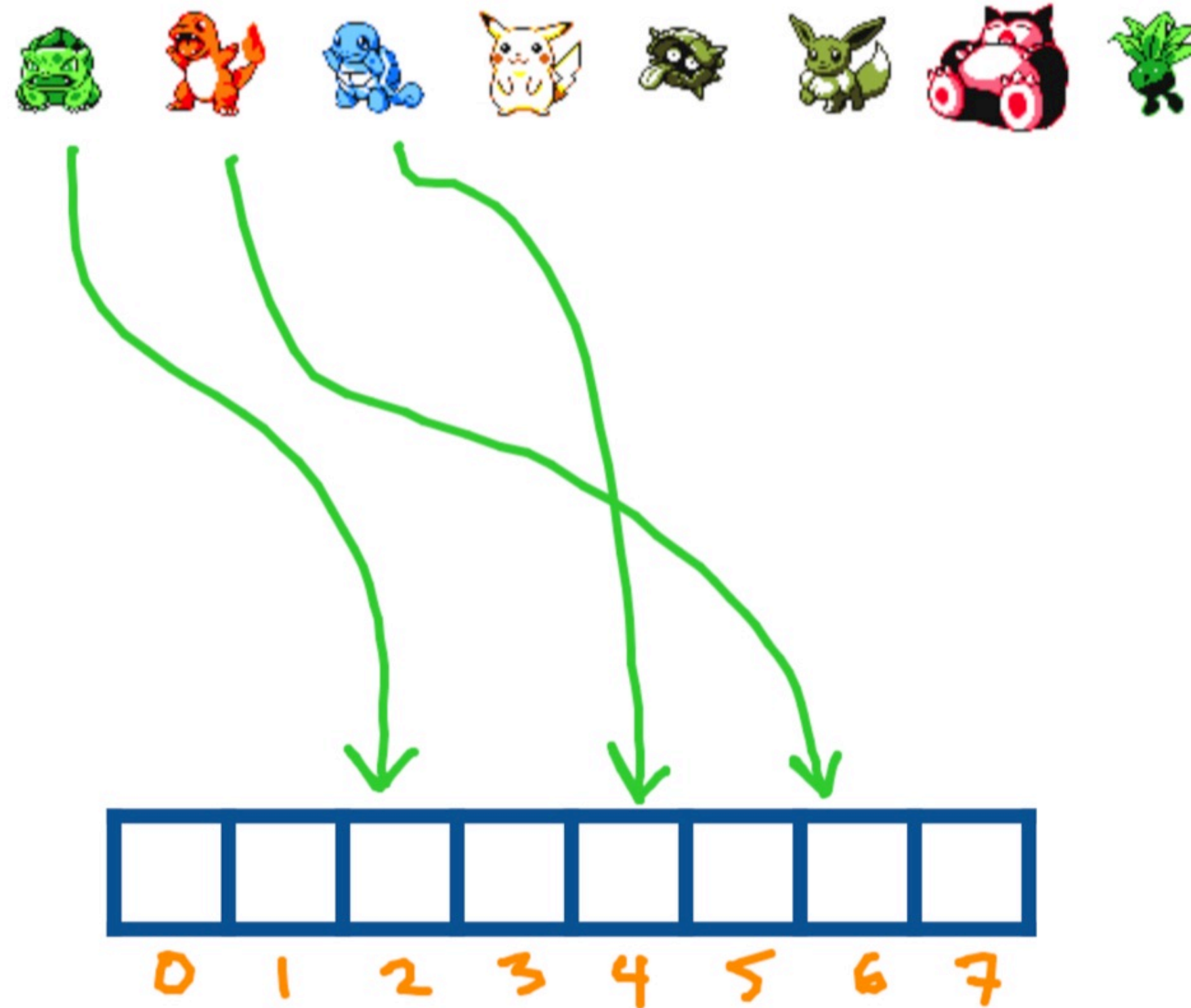
💡 Idea: use an array, but find a way to map items to array indices.

Introducing *hash functions*: map items (keys) to an array index.

Imagine our keys are **integers** (we'll come back to custom types like **Point** later).

Example: storing student ID to name.

- Green onion: 00837194 **2**
- Fire lizard: 00833462 **6**
- Blue turtle: 00749132 **4**
- Electric mouse: 00678395 **3**
- Clam: 00889321 **1**
- Fox: 00765432 **0**
- Sleepy panda: 00812493 **5**
- Radish: 00754639 **7**



Idea 1: huge
array, use ID
as index

Better idea:
map ID to
index with
function, eg $ID \% 8$

Hash functions: map the universe of keys to a restricted range (e.g. the size m of an array).

What makes a good hash function?

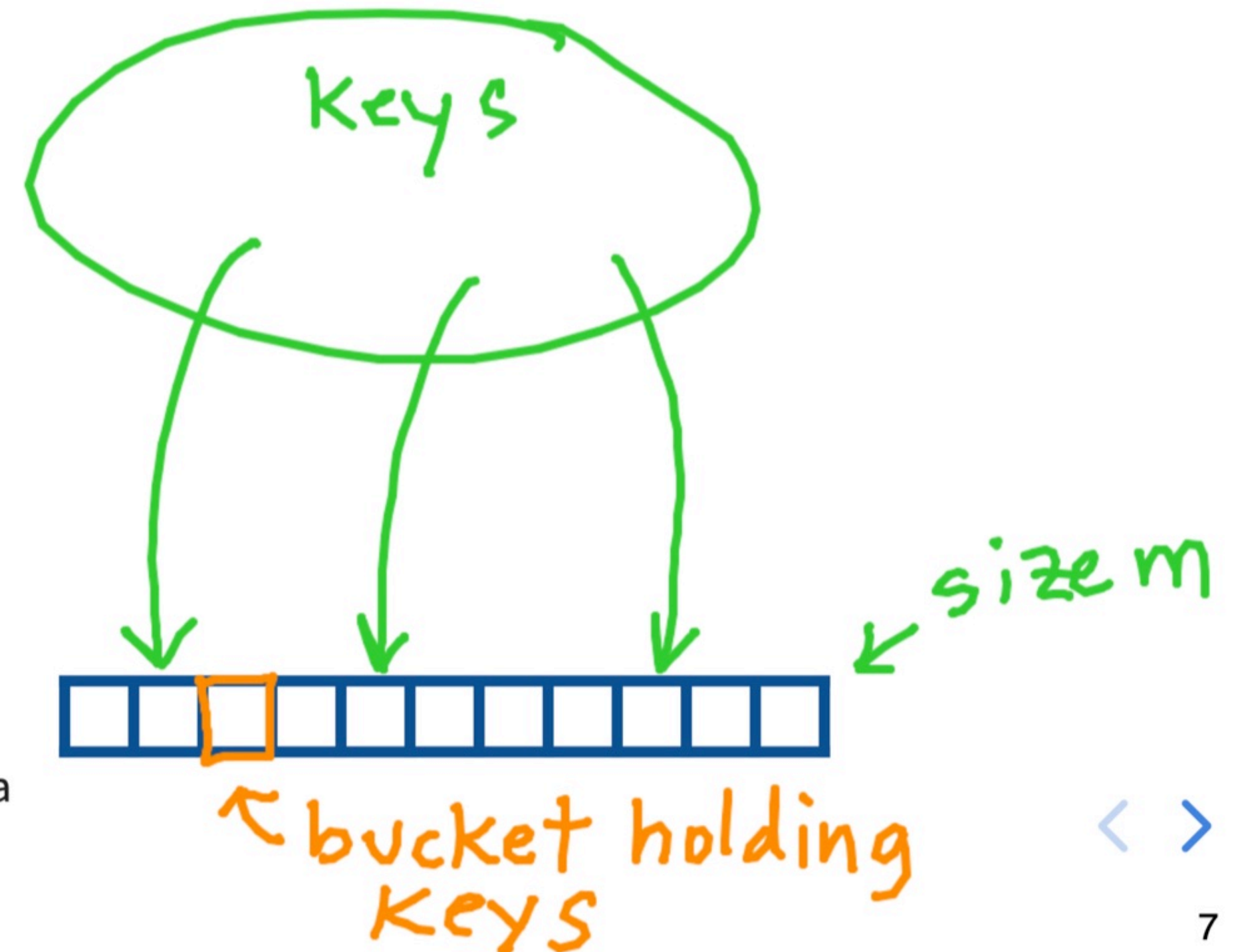
- **Deterministic:** $h(k)$ should always return the same value.
- **Fast to evaluate:** if it's expensive (e.g. $\log n$) then we don't gain anything by using a table.

The challenge is that we generally don't know the distribution of the keys.

Example: division method ($h(k) = k \bmod m$).

m	k	$h(k)$
11	25	3
11	1	1
11	17	6
13	133	3
13	7	7
13	25	12

Good rule of thumb (for division method) is that m is a prime number not too close to a power of 2.



Concern: what if multiple keys map to the same array index?

COLLISION

m	k	$h(k)$
11	25	3
11	1	1
11	17	6
11	34	1

→ map to same index



Implementing a hash table for custom types.

- Override `public int hashCode()`.
- To `add(K key)` or `put(K key, V value)`:
 1. Evaluate `key.hashCode()` to get an integer.
 2. Use the result to determine the array index.
 3. Place item in bucket at array index.
- To `get(K key)`:
 1. Evaluate `key.hashCode()` to get an integer.
 2. Use the result to determine the array index.
 3. Retrieve item in bucket at array index.
- But as we said, there may be collisions (or multiple keys that map to the same bucket). We need to iterate through items in the bucket to find the item that is associated with the `key`:
 - So we need to find the item with a `key` (i.e. check `key` equivalence).
 - Requires overriding `public boolean equals(Object otherObj)`.
- **Load factor:** $\alpha = n/m$.
 - n : number of items (`size`).
 - m : number of buckets (length of table array, i.e. capacity).
 - Higher α : potentially go through many items in one bucket to search.
 - Low α means wasted memory.
- Resize the table if load factor is outside acceptable range.
 - And be sure to rehash items (recompute table indices).



```
1 class Point {
2     public int x;
3     public int y;
4     public Point(int x, int y) {
5         this.x = x;
6         this.y = y;
7     }
8
9     @Override
10    public int hashCode() {
11        // TODO how should we hash two integer values?
12    }
13
14    @Override
15    public boolean equals(Object otherObject) {
16        Point otherPoint = (Point) otherObject;
17        return (otherPoint.x == x) && (otherPoint.y == y);
18    }
19 }
```

A detailed look into how **Java** does this (**OpenJDK**).

<https://github.com/openjdk/jdk/blob/7b0f273e37625461baa333a3ef20fbbd93647243/src/java.base/share/classes/java/util/HashMap.java#L320>

XOR

01101
11001

10100

AND

00110
00111

00110

```
1 /**
2  * Computes key.hashCode() and spreads (XORs) higher bits of hash
3  * to lower. Because the table uses power-of-two masking, sets of
4  * hashes that vary only in bits above the current mask will
5  * always collide. (Among known examples are sets of Float keys
6  * holding consecutive whole numbers in small tables.) So we
7  * apply a transform that spreads the impact of higher bits
8  * downward. There is a tradeoff between speed, utility, and
9  * quality of bit-spreading. Because many common sets of hashes
10 * are already reasonably distributed (so don't benefit from
11 * spreading), and because we use trees to handle large sets of
12 * collisions in bins, we just XOR some shifted bits in the
13 * cheapest possible way to reduce systematic lossage, as well as
14 * to incorporate impact of the highest bits that would otherwise
15 * never be used in index calculations because of table bounds.
16 */
17 static final int hash(Object key) {
18     int h;
19     return (key == null) ? 0 : (h = key.hashCode()) ^ (h >>> 16);
20 }
```

- Table (array) size is always a power of 2.
- Table index computed from **hash(key) & (m - 1)** where *m* is the capacity (table length).
- ^ means bitwise XOR (exclusive OR): e.g. **01101 ^ 11001** is
- >>> means unsigned right shift (here, by 16 bits): e.g. **01101010 >>> 4** is **00000110**
- & means bitwise AND: e.g. **00110 & 111** is
- In your **Terminal**, type **jshell** and use **Integer.toBinaryString** to try these out! **Ctrl-D** to exit.



What exactly is **hash(key) & (m - 1)** doing?

Remember that **Java** (**OpenJDK**) picks the table size to be a power of 2.

- Example: $m = 16$ is **10000**, so $m - 1$ is **1111**

- What is **15 & 15**? **1111 & 1111 = 1111 = 15**

- What is **73 & 15**?

```
      9
    00011111
  00010001
  -----
    00010001
      9
```


Exercise: compute the hash table indices for the following **Points** (using **Java**'s technique).

```

1 static final int hash(Object key) {
2     int h;
3     return (key == null) ?
4         0 : (h = key.hashCode()) ^ (h >>> 16);
5 }

```

- Starting with a table size $m = 16$.
- Table index = $\text{hash}(\text{key}) \& (m - 1)$.
- Point 1: (1, 1)
- Point 2: (12345, 678)

```

1 class Point {
2     public int x;
3     public int y;
4     public Point(int x, int y) {
5         this.x = x;
6         this.y = y;
7     }
8
9     @Override
10    public int hashCode() {
11        return 31 * x + y;
12    }
13
14    @Override
15    public boolean equals(Object otherObject) {
16        Point otherPoint = (Point) otherObject;
17        return (otherPoint.x == x) && (otherPoint.y == y);
18    }
19 }

```

Point	hashcode h $31 * x + y$	rightshift $h \gg \gg 16$	$h \wedge (h \gg \gg 16)$ final hash	table index $h \& 15$
(1, 1)	32	0	32	0 ... 0000 & 1111
(12345, 678)	383373	5	383368	8 ... 1000 & 1111
(70000, 42)	2170042	33	2170011	11 ... 1011 & 1111

Exercise: compute the hash table indices for the following **Points** (using **Java**'s technique).

```
1 static final int hash(Object key) {  
2     int h;  
3     return (key == null) ?  
4         0 : (h = key.hashCode()) ^ (h >>> 16);  
5 }
```

a quick way to compute mod

- Starting with a table size $m = 16$.
- Table index = $\text{hash}(\text{key}) \ \& \ (m - 1)$. \leftarrow
- Point 1: (1, 1)
- Point 2: (12345, 678)

```
1 class Point {  
2     public int x;  
3     public int y;  
4     public Point(int x, int y) {  
5         this.x = x;  
6         this.y = y;  
7     }  
8  
9     @Override  
10    public int hashCode() {  
11        return 31 * x + y;  
12    }  
13  
14    @Override  
15    public boolean equals(Object otherObject) {  
16        Point otherPoint = (Point) otherObject;  
17        return (otherPoint.x == x) && (otherPoint.y == y);  
18    }  
19 }
```

- HashMap table size is power of 2

- hash is designed to spread keys across buckets

$\ggg 16$ is taking higher bits

$h \wedge (h \ggg 16)$ is mixing high and low order bits to better spread in table

Additional notes:

- For more hash functions, see: https://en.wikipedia.org/wiki/List_of_hash_functions
- **Lab 7** due tonight.
- **Homework 8** due on Thursday 4/24: use a **TreeMap** to solve a problem and then implement a DIY-version based on what your algorithm needs.
- **Next class:** how can we handle collisions?



Always acknowledge
any use of
generative AI
in code

Two large, red, 3D hash symbols (#) are positioned to the left of the handwritten text, one above the other.